

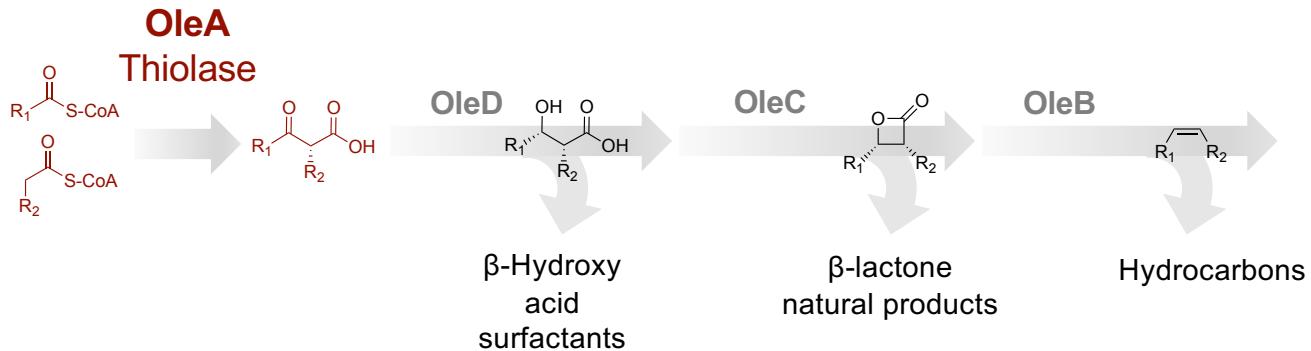
Machine learning-guided prediction of OleA thiolase activity and substrate scope

Serina L. Robinson,^{1,2,3} Megan D. Smith,^{2,3} Jack E. Richman,³ Kelly G. Aukema,³ Lawrence P. Wackett³

¹Graduate Program in Bioinformatics and Computational Biology, University of Minnesota 111 S. Broadway, Suite 300, Rochester, MN, 55904, USA; ²Graduate Program in Microbiology, Immunology, and Cancer Biology, University of Minnesota, 689 23rd Ave SE, Minneapolis, MN, 55455, USA; ³BioTechnology Institute, University of Minnesota, 1479 Gortner Avenue, Saint Paul, MN, 55108, USA

Main text figures

A



B

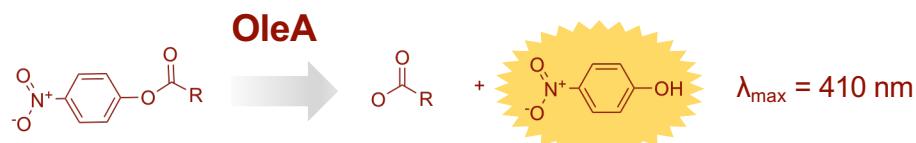


Figure 1. Thiolase enzymes in the OleA family catalyze a head-to-head Claisen condensation of two acyl-CoA substrates (maroon) as the first committed step in production of value-added metabolites such as surfactants, pharmaceuticals and hydrocarbons. R₁, R₂ in native OleA pathways: C₈ – C₁₆. (B) OleA reacts with various *para*-nitrophenyl esters to produce the corresponding carboxylic acids and *para*-nitrophenol chromophore, providing a rapid readout for enzyme activity.

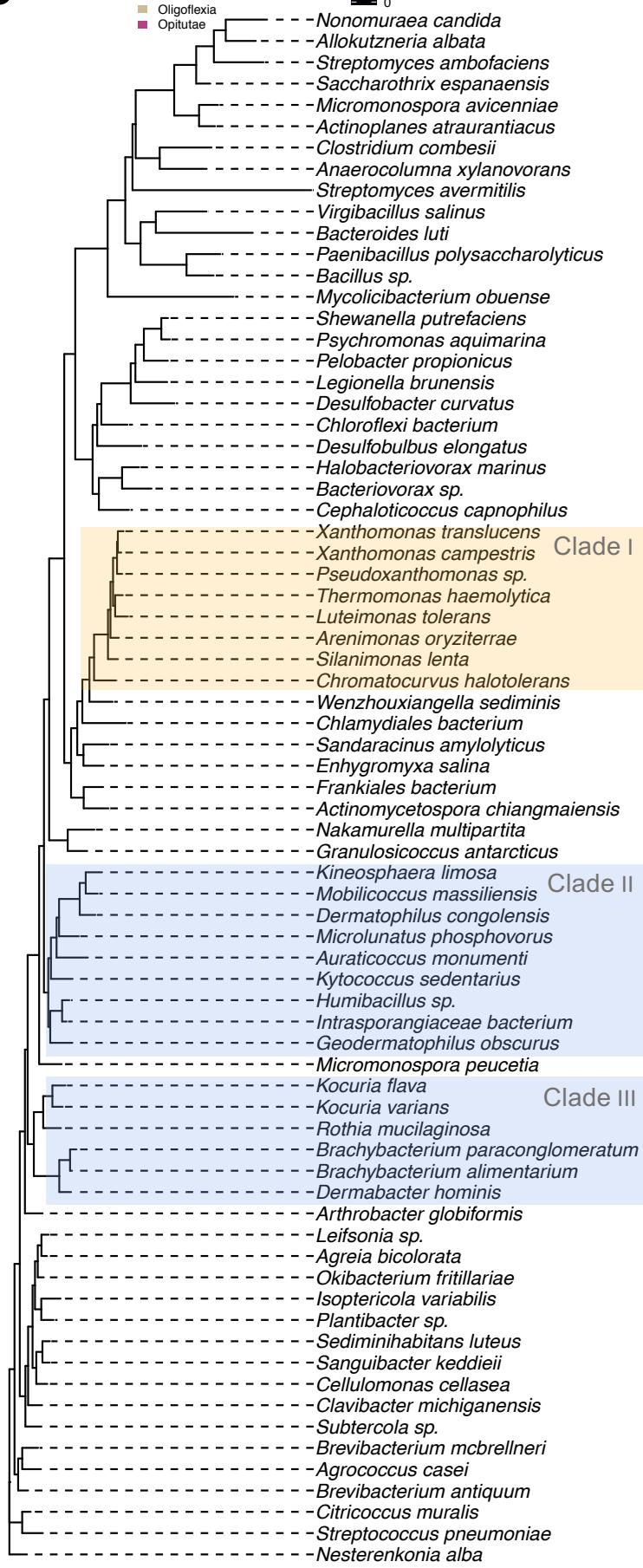
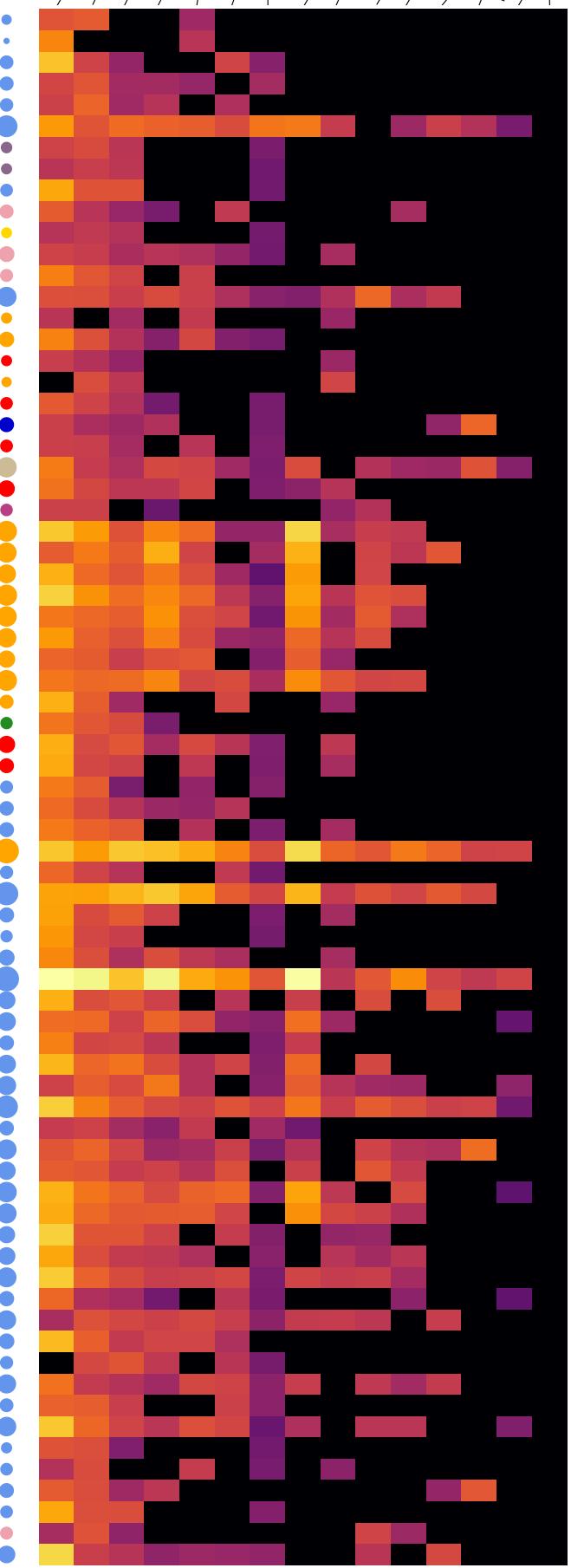
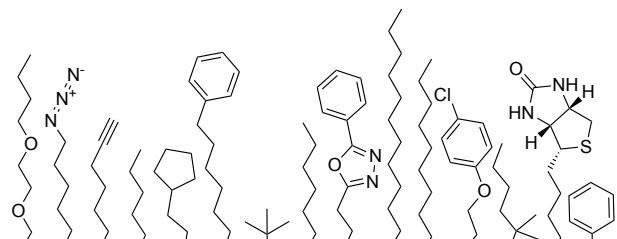
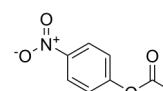
Average enzyme activity*

Taxonomic class

Enzyme activity per substrate†

- 0.5
- 1.0
- 1.5
- 2.0

- Actinobacteria
- Bacilli
- Bacteroidia
- Chlamydiae
- Clostridia
- Delta-proteobacteria
- Gamma-proteobacteria
- Green non-sulfur bacteria
- Oligoflexia
- Opitutae



Clade I

Clade II

Clade III

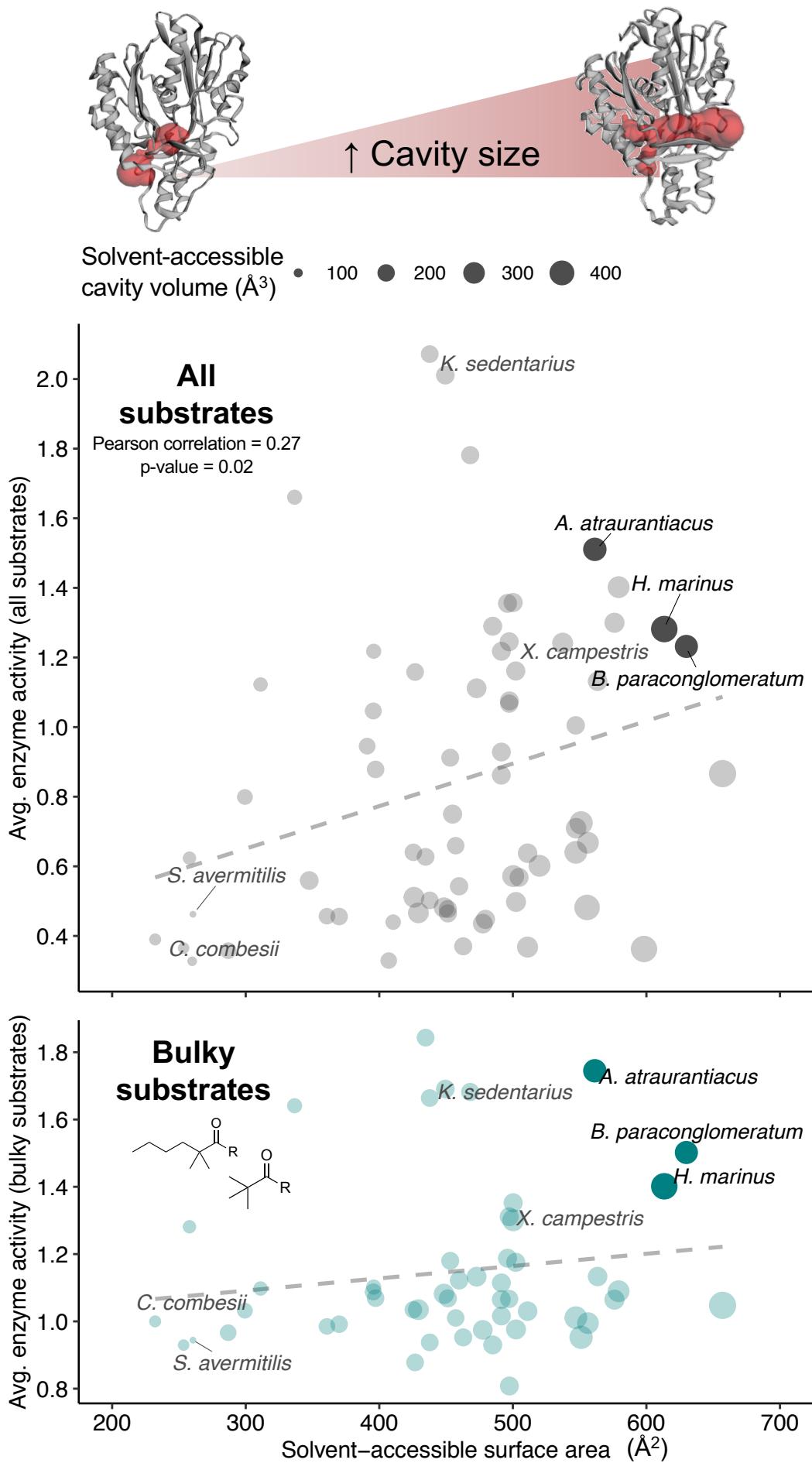


Figure 3. (A) Predicted solvent-accessible cavity volumes of the 73 OleA enzymes analyzed in this study ranged from 95.3 to 478.9 \AA^3 (B) Relationship between the solvent-accessible surface area with the average enzymatic activity across all substrates (black) and 'bulky' substrates (teal) with tertiary α -carbons (*p*NP trimethylacetate and *p*NP dimethyl hexanoate). *A. atraurantiacus*, *B. paraconglomeratum*, and *H. marinus* enzymes have among the highest calculated cavity volumes and the highest preferences for substrates with tertiary α -carbons.

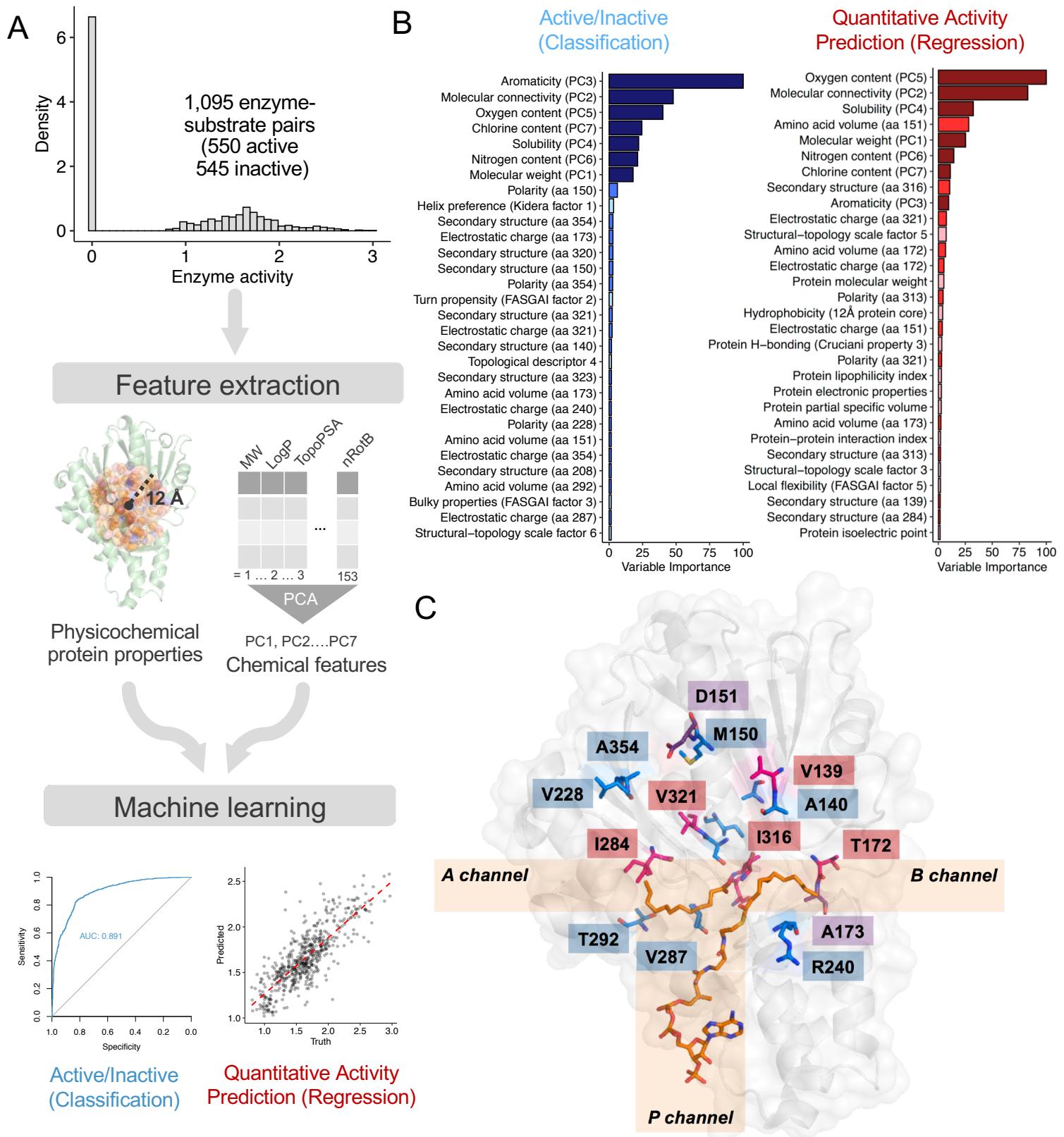


Figure 4. (A) Machine learning workflow (B) Variable importance scores for classification and regression models (C) Important residues mapped onto the *X. campestris* 4KU5 structure with fatty acid substrates bound. Residue colors correspond to variable importance in the classification model (blue), regression model (red), or both models (purple).

Table 1. Machine learning classification results. [†]CI = Confidence interval. See Fig S3 for extended results.

| Machine learning algorithm | Training classification accuracy | Testing classification accuracy | Testing 95% CI [†] |
|----------------------------|----------------------------------|---------------------------------|-----------------------------|
| Random Forest | 0.826 | 0.839 | (0.789, 0.880) |
| Feedforward Neural Network | 0.732 | 0.777 | (0.722, 0.826) |
| Naïve Bayes | 0.586 | 0.645 | (0.585, 0.701) |

Table 2. Machine learning regression results. See Fig S3 for extended results.

| Machine learning algorithm | Training RMSE | Training R ² | Testing RMSE | Testing R ² |
|--|---------------|-------------------------|--------------|------------------------|
| Random Forest | 0.254 | 0.625 | 0.219 | 0.745 |
| Multivariate adaptive regression splines | 0.276 | 0.557 | 0.252 | 0.642 |
| Elastic Net | 0.275 | 0.549 | 0.278 | 0.564 |

Supplemental Figures

Table S1. Accession numbers, organisms, average activity across all substrates, and total number of substrates accepted (out of 15). [†]Enzyme activity is measured in the log₁₀ of nmol pNP produced over the course of one hour by 200 (**uL**) *E. coli* BL21 culture with an OD of 1.0 heterologously expressing OleA averaged across three independent biological replicates and across all substrates.

| NCBI Accession | Organism | Average enzyme activity [†] | Total number of substrates accepted |
|----------------|--|--------------------------------------|-------------------------------------|
| WP_012802190 | <i>Kytococcus sedentarius</i> | 2.072 | 14 |
| WP_088919872 | <i>Granulosicoccus antarcticus</i> | 2.011 | 14 |
| WP_040156198 | <i>Mobilicoccus massiliensis</i> | 1.781 | 13 |
| WP_068469535 | <i>Kocuria varians</i> strain G6 | 1.66 | 14 |
| WP_097326105 | <i>Actinoplanes atrauranticus</i> | 1.511 | 13 |
| WP_114959563 | <i>Thermomonas haemolytica</i> | 1.402 | 11 |
| WP_117315975 | <i>Chromatocurvus halotolerans</i> | 1.358 | 11 |
| WP_003480168 | <i>Xanthomonas translucens</i> pv. <i>graminis</i> | 1.355 | 11 |
| WP_101686474 | <i>Dermabacter hominis</i> | 1.3 | 11 |
| WP_076585679 | <i>Luteimonas tolerans</i> | 1.29 | 11 |
| WP_096909276 | <i>Halobacteriovorax marinus</i> strain BE01 | 1.282 | 13 |
| NP_635607 | <i>Xanthomonas campestris</i> | 1.245 | 10 |
| WP_003805730 | <i>Arthrobacter globiformis</i> | 1.241 | 10 |
| WP_126985334 | <i>Brachybacterium paraconglomeratum</i> | 1.232 | 12 |
| WP_079726385 | <i>Oikibacterium fritillariae</i> | 1.218 | 11 |
| WP_082133371 | <i>Mycolicibacterium obuense</i> | 1.218 | 12 |
| WP_022969495 | <i>Arenimonas oryziterrae</i> | 1.161 | 10 |
| WP_104244666 | <i>Subtercola</i> sp. Z2020 | 1.158 | 11 |
| WP_086472704 | <i>Plantibacter</i> sp. VKM Ac-1784 | 1.132 | 11 |
| WP_058859095 | <i>Kocuria flava</i> | 1.122 | 11 |
| WP_034630949 | <i>Cellulomonas cellasea</i> | 1.111 | 11 |
| WP_130552392 | <i>Pseudoxanthomonas</i> sp. NML171200 | 1.075 | 9 |
| WP_091632192 | <i>Micromonospora peucetia</i> | 1.067 | 9 |
| WP_026862459 | <i>Intrasporangiaceae</i> bacterium URHB0013 | 1.046 | 10 |
| WP_096166078 | <i>Brachybacterium alimentarium</i> | 1.005 | 9 |
| WP_070193933 | <i>Humibacillus</i> sp. DSM 29435 | 0.945 | 8 |
| WP_044441590 | <i>Agreia bicolorata</i> | 0.928 | 9 |
| WP_033420292 | <i>Nesterenkonia alba</i> | 0.912 | 9 |
| WP_028771085 | <i>Silanimonas lenta</i> | 0.878 | 8 |
| WP_053234550 | <i>Sandaracinus amyloyticus</i> | 0.866 | 8 |
| WP_056164922 | <i>Leifsonia</i> sp. Leaf325 | 0.862 | 8 |
| WP_021267235 | <i>Bacteriovorax</i> sp. BAL6 X | 0.799 | 8 |
| WP_090596165 | <i>Auraticoccus monumenti</i> | 0.75 | 7 |
| WP_090922125 | <i>Paenibacillus polysaccharolyticus</i> | 0.725 | 8 |
| WP_100422947 | <i>Sediminibhabitans luteus</i> | 0.709 | 6 |
| WP_028865527 | <i>Psychromonas aquimarna</i> | 0.678 | 7 |
| WP_013837355 | <i>Isoptericola variabilis</i> | 0.668 | 8 |
| WP_028327376 | <i>Dermatophilus congolensis</i> | 0.659 | 6 |
| OGN97459 | <i>Chloroflexi</i> bacterium RBG 13 51 36 | 0.648 | 7 |
| WP_015749354 | <i>Nakamurella multipartita</i> | 0.64 | 6 |
| WP_012948569 | <i>Geodermatophilus obscurus</i> | 0.64 | 6 |
| WP_106091803 | <i>Enhygromyxa salina</i> strain SWB007 | 0.638 | 6 |
| WP_101620581 | <i>Brevibacterium antiquum</i> | 0.627 | 6 |
| WP_005508112 | <i>Rothia mucilaginosa</i> | 0.623 | 7 |
| WP_026204600 | <i>Actinomycetospora chiangmaiensis</i> | 0.601 | 6 |
| WP_015101236 | <i>Saccharothrix espanaeensis</i> | 0.572 | 6 |
| WP_116651321 | <i>Wenzhouxiangella sediminis</i> | 0.568 | 5 |
| WP_092493676 | <i>Virgibacillus salinus</i> | 0.559 | 6 |
| WP_015490789 | <i>Clavibacter michiganensis</i> | 0.543 | 5 |
| WP_053141838 | <i>Streptomyces ambofaciens</i> | 0.542 | 5 |
| WP_076469753 | <i>Micromonospora avicenniae</i> | 0.511 | 5 |
| WP_006592558 | <i>Kineosphaera limosa</i> | 0.502 | 5 |
| WP_012868615 | <i>Sanguibacter keddieii</i> | 0.498 | 5 |
| WP_119951226 | <i>Frankiales</i> bacterium YIM 75000 | 0.482 | 5 |
| WP_095301808 | <i>Bacillus</i> sp. 7586-K | 0.482 | 4 |
| CVN04163 | <i>Streptococcus pneumoniae</i> | 0.477 | 5 |
| WP_051554118 | <i>Desulfobulbus elongatus</i> | 0.466 | 5 |
| WP_115933252 | <i>Citricoccus muralis</i> | 0.465 | 4 |
| WP_010983715 | <i>Streptomyces avermitilis</i> | 0.462 | 4 |
| WP_020589018 | <i>Desulfobacter curvatus</i> | 0.457 | 5 |
| WP_086990725 | <i>Agrococcus casei</i> | 0.456 | 5 |
| WP_068712879 | <i>Cephalotilicoccus capnophilus</i> | 0.448 | 5 |
| WP_088206054 | <i>Chlamydiales</i> bacterium SCGC AG-110-P3 | 0.44 | 4 |
| WP_013865923 | <i>Microlunatus phosphovorus</i> | 0.435 | 4 |
| WP_099839552 | <i>Clostridium combesi</i> | 0.39 | 4 |
| WP_005885358 | <i>Brevibacterium mcbrellneri</i> | 0.37 | 4 |
| WP_061782733 | <i>Shewanella putrefaciens</i> | 0.368 | 4 |
| WP_084558798 | <i>Anaerocolumna xylanovorans</i> | 0.365 | 4 |
| ABL00697 | <i>Pelobacter propionicus</i> DSM 2379 | 0.362 | 4 |
| WP_073399371 | <i>Bacteroides luti</i> | 0.358 | 4 |
| WP_063765283 | <i>Nonomuraea candida</i> | 0.329 | 3 |
| WP_058441655 | <i>Legionella brunensis</i> | 0.327 | 3 |
| WP_030433686 | <i>Allokutzneria albata</i> | 0.243 | 2 |

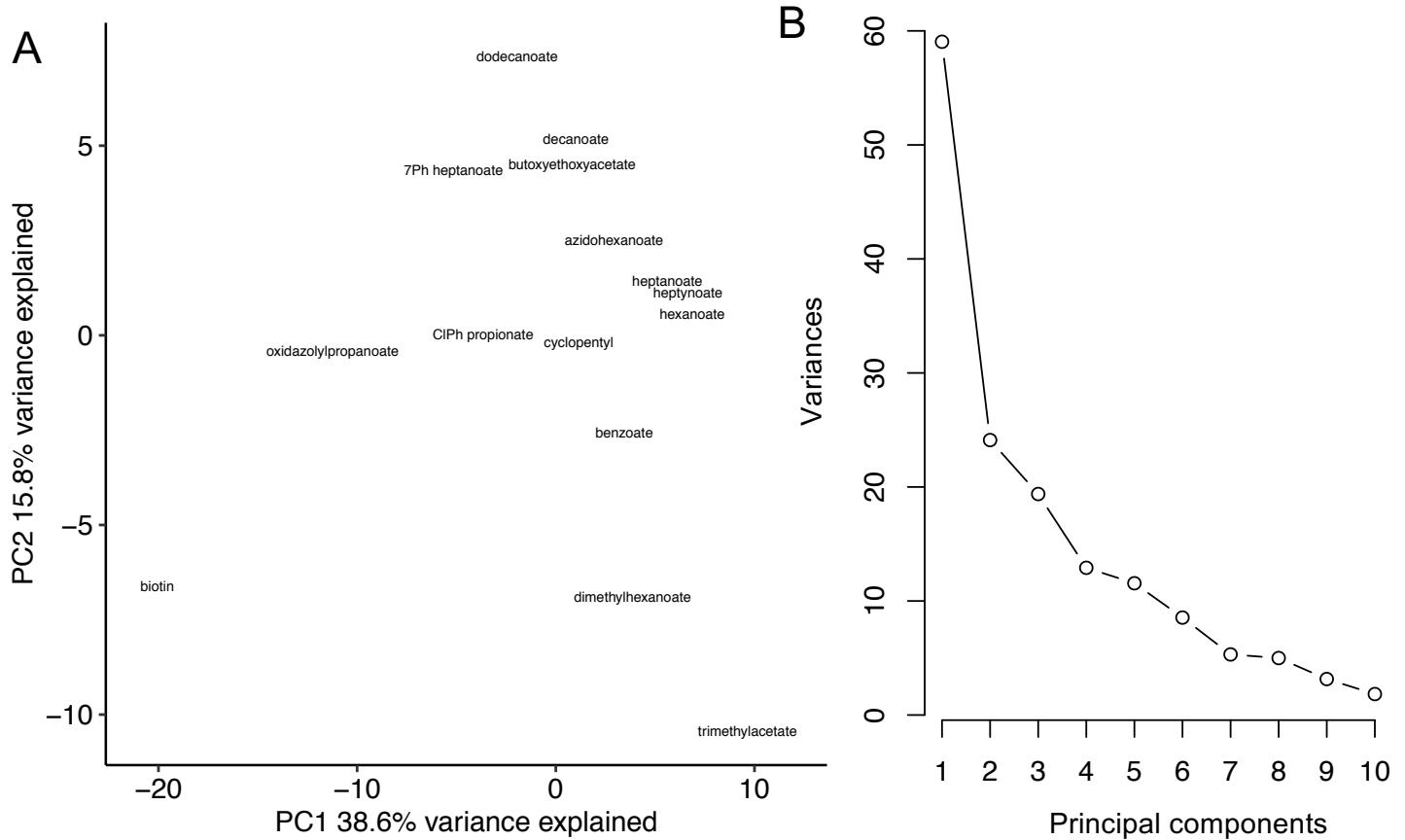


Figure S2. (A) Dimensionality reduction of molecular descriptors of 15 pNP substrates using principal components analysis. Substrates are plotted by their first two principal components that cumulatively explain up to 54.4% of the variance between compounds (B) Screeplot of principal components reveals 'elbow' after seventh principal component. (C) Absolute values of loadings of 149 molecular descriptors into 7 principal components roughly corresponding to molecular weight, molecular connectivity, aromaticity, solubility, and O, N and Cl content.

A Classification Results

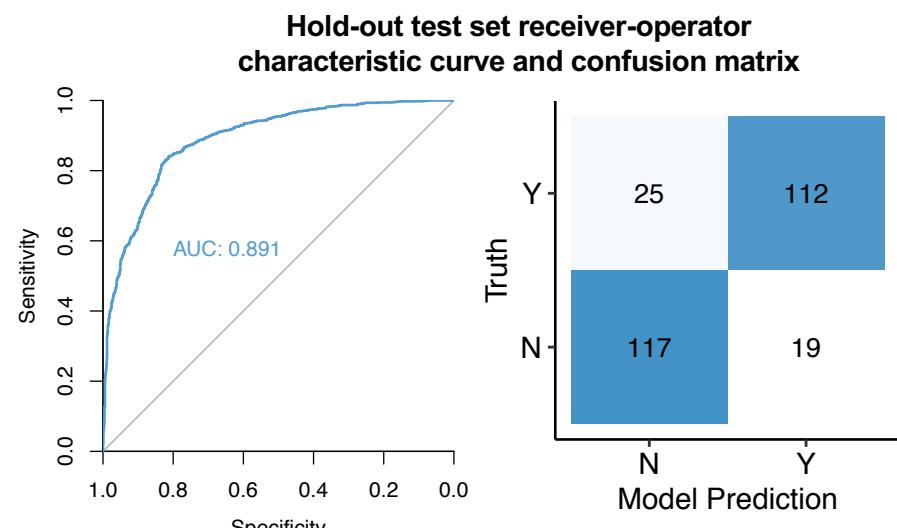
Model comparison and hyperparameter tuning (10-fold cross-validation)

| Machine learning algorithm | Training classification accuracy | Training Cohen's kappa | Testing classification accuracy | Testing 95% confidence interval | Best model hyperparameters |
|----------------------------|----------------------------------|------------------------|---------------------------------|---------------------------------|---|
| Random Forest | 0.826 | 0.651 | 0.839 | (0.789, 0.880) | mtry = 334, splitrule = extratrees, min.node.size = 1 |
| Feedforward Neural Network | 0.732 | 0.170 | 0.777 | (0.722, 0.826) | size = 3, decay = 0.5 |
| Naïve Bayes | 0.586 | 0.645 | 0.645 | (0.585, 0.701) | fL = 0, adjust = 1, useKernel = F |

Average classification results 10 different train-test splits

Training set accuracy 0.819 ± 0.009

Testing set accuracy 0.808 ± 0.017



B Regression Results

Model comparison and tuning (10-fold cross-validation)

| Machine learning algorithm | Training root-mean-square error (RMSE) | Training R ² | Training mean absolute error | Testing RMSE | Testing R ² | Best model hyperparameters |
|--|--|-------------------------|------------------------------|--------------|------------------------|---|
| Random Forest | 0.254 | 0.625 | 0.198 | 0.219 | 0.745 | mtry = 168, splitrule = variance, min.node.size = 5 |
| Multivariate adaptive regression splines | 0.276 | 0.557 | 0.217 | 0.252 | 0.642 | nprune = 20, degree = 1 |
| Elastic Net | 0.275 | 0.549 | 0.214 | 0.278 | 0.564 | alpha = 1 (Lasso) lambda = 0.0037 |

Average regression results 10 different training-test splits

Training set RMSE 0.249 ± 0.007

Testing set RMSE 0.245 ± 0.019

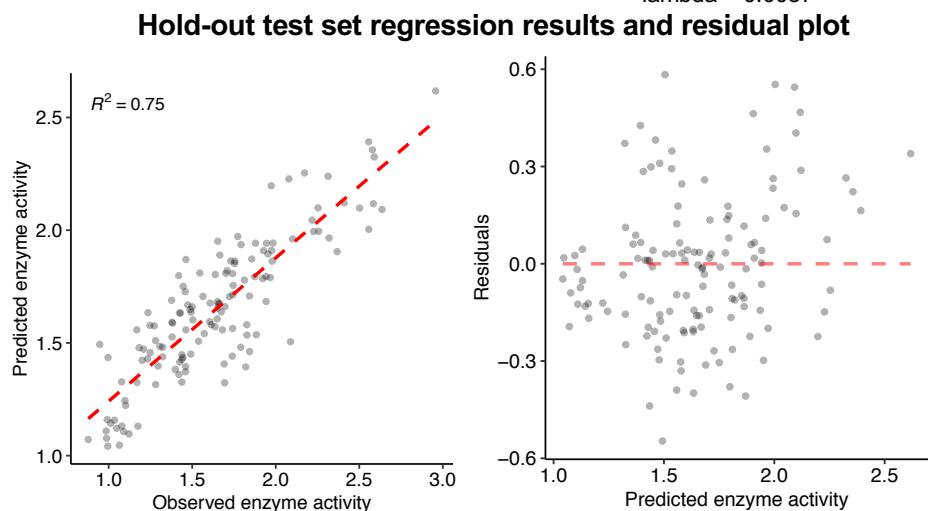
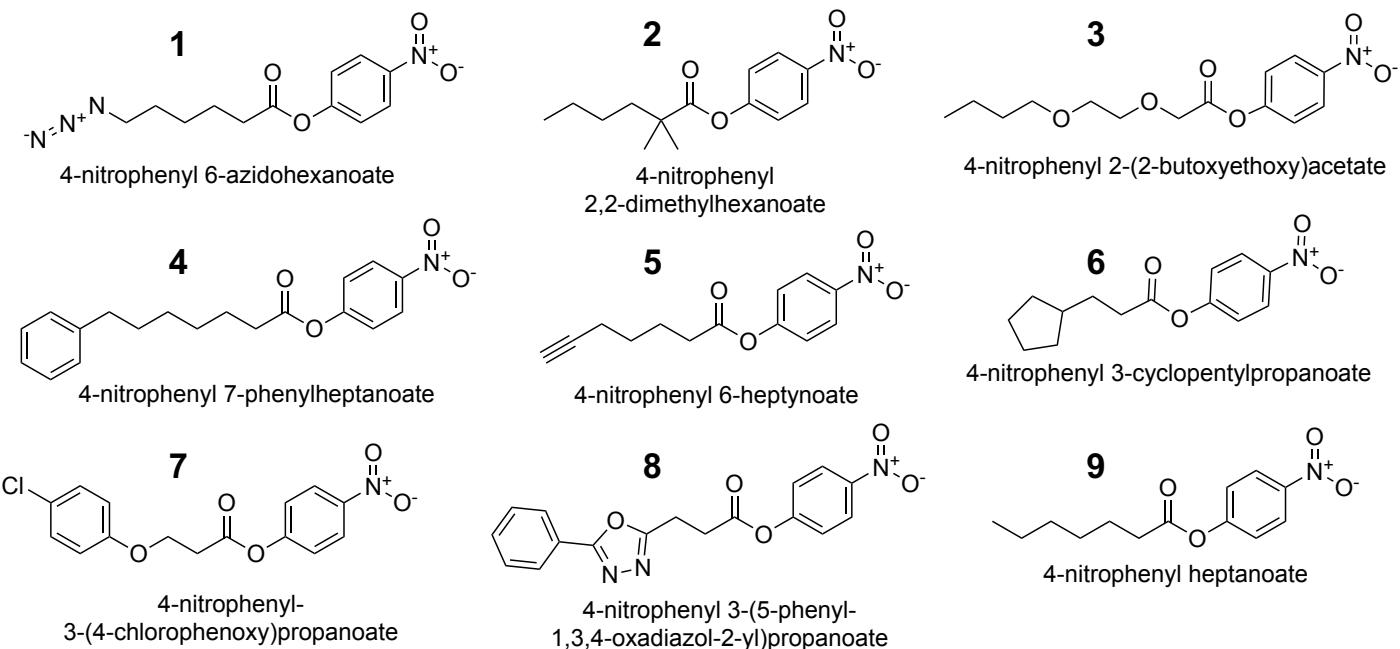
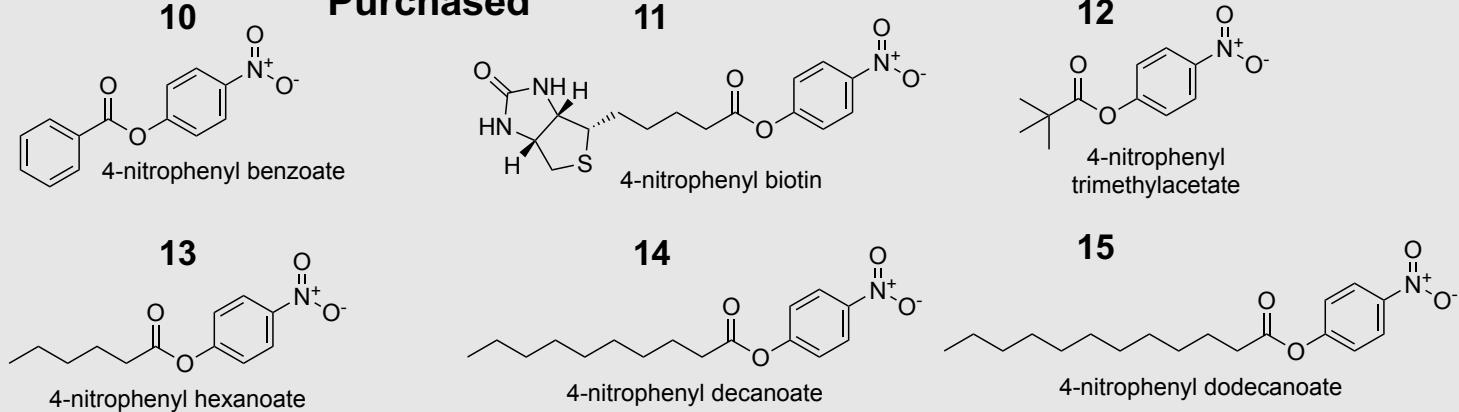


Figure S3. (A) Machine learning classification results and (B) regression results.

A

Synthesized**Purchased**

B

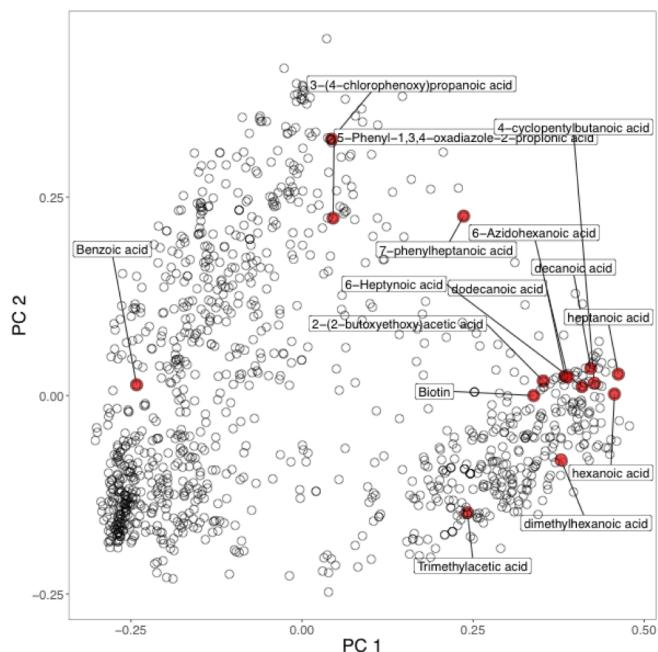
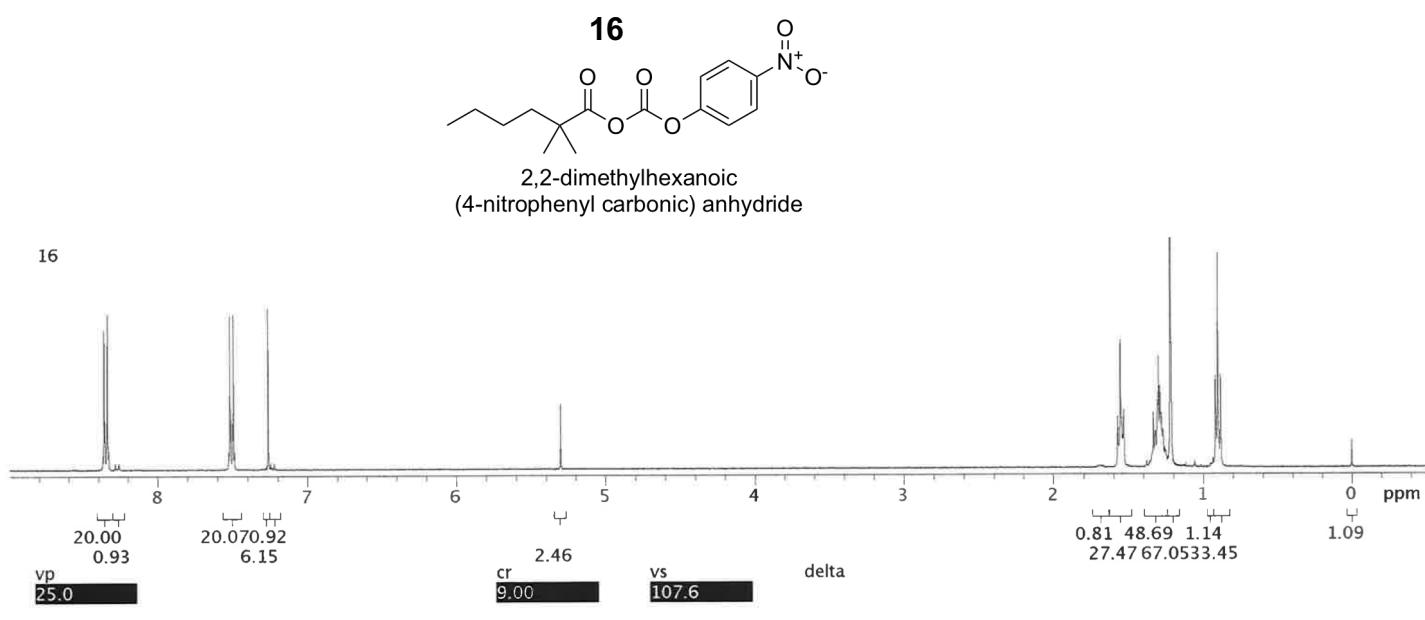
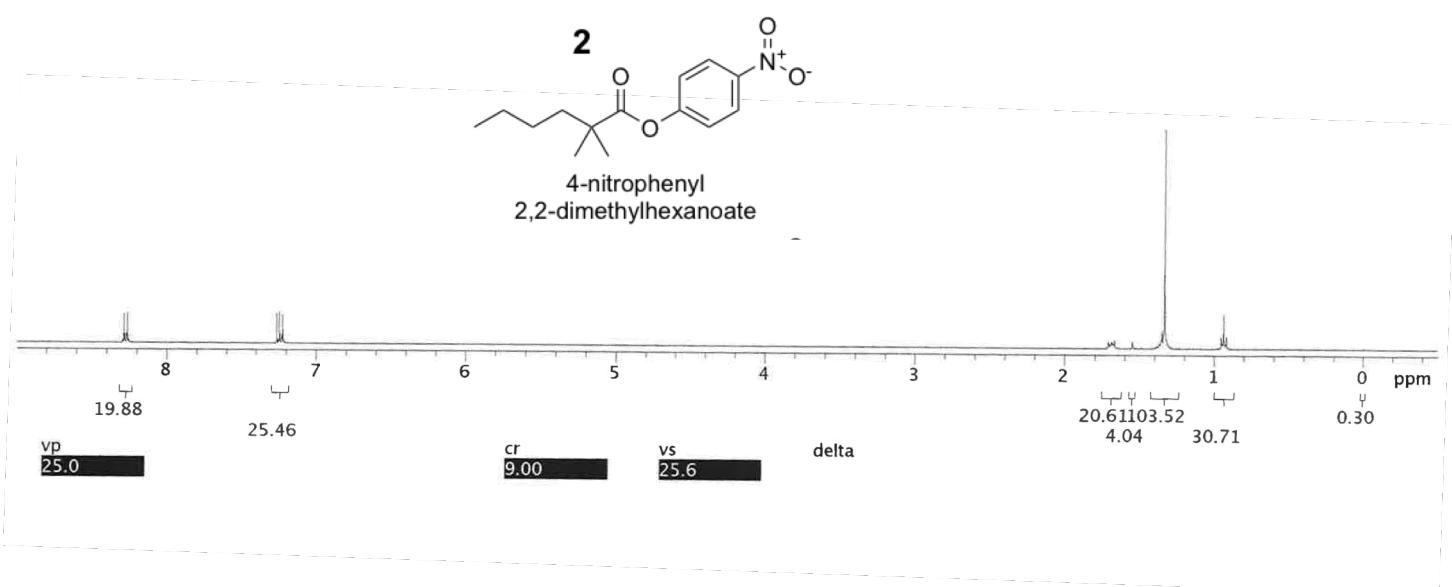
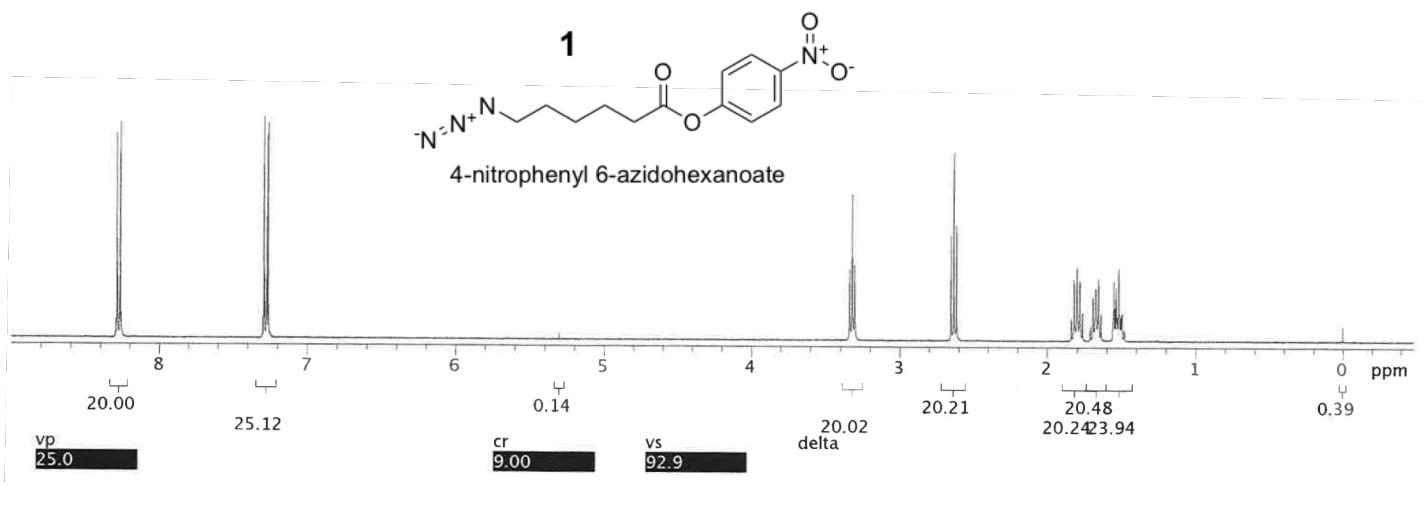
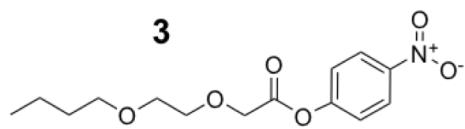
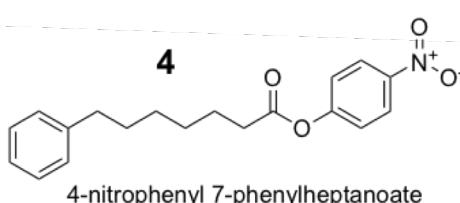
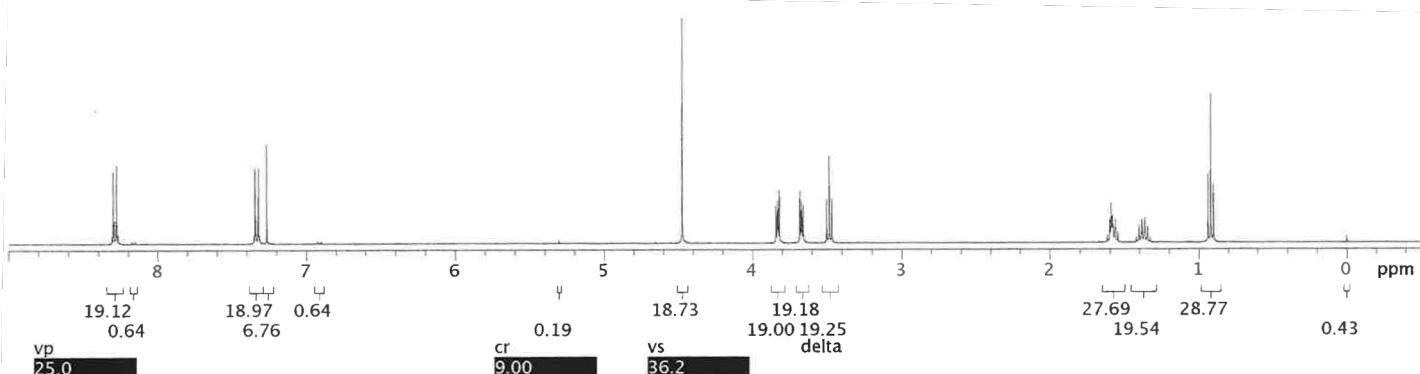


Figure S4. (1) Chemical structures of *p*NP esters synthesized and purchased in this study. (B) Tanimoto clustering of 15 *p*NP substrates (maroon) compared to the sequence space of commercially-available carboxylic acid substrates from Sigma-Aldrich.

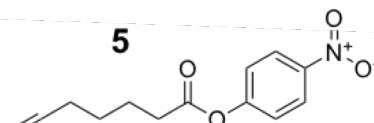
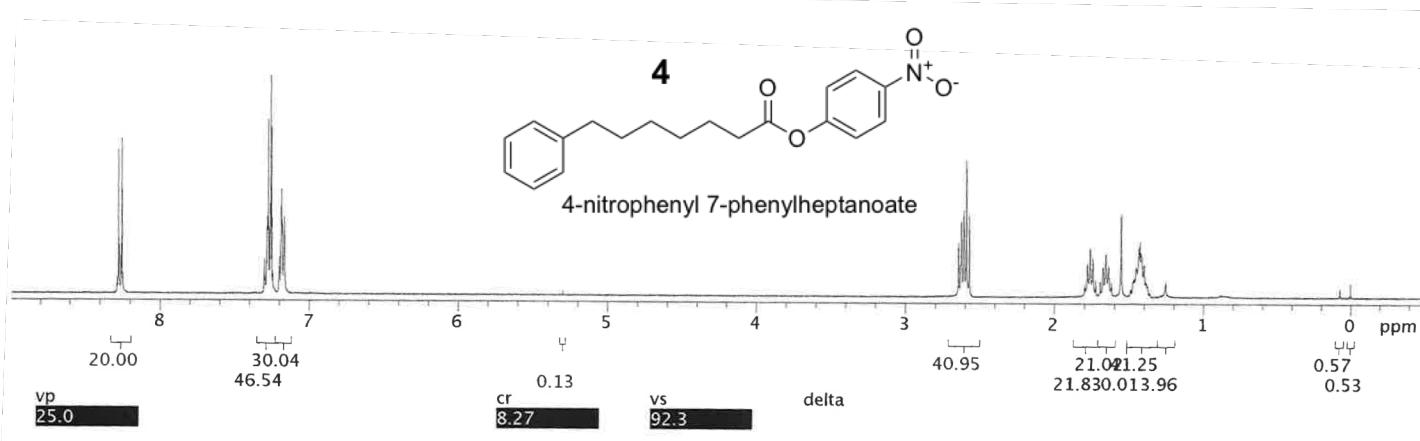




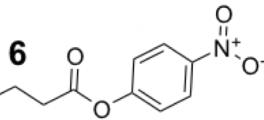
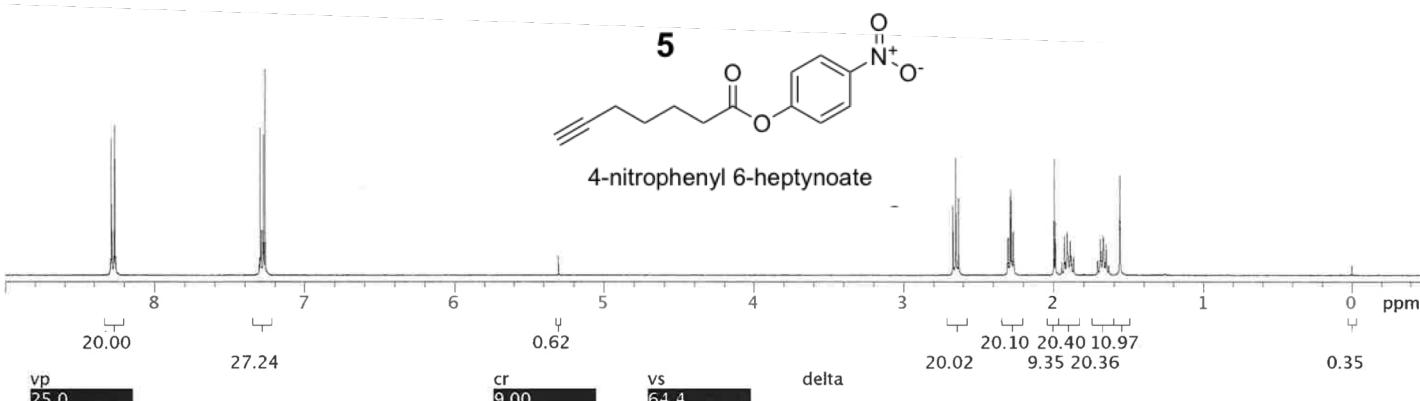
4-nitrophenyl 2-(2-butoxyethoxy)acetate



4-nitrophenyl 7-phenylheptanoate



4-nitrophenyl 6-heptynoate



4-nitrophenyl 3-cyclopentylpropanoate

